

EXHIBIT 125

11/21/2024

Richard Kadrey, et al. v. Meta Platforms, Inc. Yann LeCun, Ph.D.
Highly Confidential - Attorneys' Eyes Only - Under the Protective Order

Page 1

IN THE UNITED STATES DISTRICT COURT
FOR THE NORTHERN DISTRICT OF CALIFORNIA
SAN FRANCISCO DIVISION

IN RE MATTER OF:)	
RICHARD KADREY, et al.,)	
Plaintiff,)	
vs.)	C.A. NO.:
META PLATFORMS, INC.,)	3:23-cv-03417-VC
Defendant.)	

VIDEOTAPED DEPOSITION OF YANN LeCUN, Ph.D.

Palo Alto, California

Thursday, November 21, 2024

** HIGHLY CONFIDENTIAL - ATTORNEYS EYES ONLY **

UNDER THE PROTECTIVE ORDER

Stenographically Reported by:

HEATHER J. BAUTISTA, CSR, CRR, RPR, CLR

Realtime Systems Administrator

California CSR License #11600

Oregon CSR License #21-0005

Washington License #21009491

Nevada CCR License #980

Texas CSR License #10725

DIGITAL EVIDENCE GROUP
1730 M. Street, NW, Suite 812
Washington, D.C. 20036
(202) 232-0646

11/21/2024

Richard Kadrey, et al. v. Meta Platforms, Inc.

Yann LeCun, Ph.D.

Highly Confidential - Attorneys' Eyes Only - Under the Protective Order

Page 24

1 If you have an understanding of the question 10:27:57

2 independent of that, by means you can provide it. 10:28:00

3 MS. GEMAN: That question was plainly not 10:28:04

4 one that would beget that admonition, so I would 10:28:05

5 respectfully request that you keep your objections 10:28:10

6 to form. 10:28:12

7 MR. WEINSTEIN: I -- I understand that. I 10:28:12

8 understand you weren't meaning to get that, but the 10:28:13

9 scope of your question is pretty broad and -- 10:28:15

10 Q. (By Ms. Geman) Do you need to hear the 10:28:18

11 question again? 10:28:19

12 A. No. I think there are two aspects. The 10:28:20

13 first one is a completely practical one which is 10:28:22

14 that those systems are -- are trained -- and the 10:28:29

15 answer is different depending on which version of -- 10:28:32

16 (Stenographer clarification.) 10:28:36

17 THE WITNESS: -- LLaMA which is the -- 10:28:37

18 Q. (By Ms. Geman) So I will -- I can make the 10:28:39

19 question more precise. 10:28:40

20 A. Yes. 10:28:41

21 Q. Do you agree with Meta's decision to fail 10:28:42

22 to disclose the training sources for LLaMA 2, 3, and 10:28:43

11/21/2024

Richard Kadrey, et al. v. Meta Platforms, Inc.

Yann LeCun, Ph.D.

Highly Confidential - Attorneys' Eyes Only - Under the Protective Order

		Page 25
1	4?	10:28:50
2	A. So 4 doesn't exist yet.	10:28:53
3	Q. I understand that.	10:28:55
4	A. So there's two things you can list; right?	10:28:56
5	So you can list the -- the set of datasets that we	10:29:01
6	used to train the system which was done for LLaMA 1.	10:29:10
7	It was not done for LLaMA 2 and 3.	10:29:13
8	You can also just redistribute the data and	10:29:16
9	that wasn't done for any of the models. For LLaMA	10:29:19
10	1, all those sources were publicly downloadable from	10:29:24
11	the open Internet.	10:29:28
12	For LLaMA 3, that wouldn't be possible	10:29:30
13	because some of the data comes from -- it's internal	10:29:32
14	data essentially from users of Meta's services, and	10:29:37
15	so it cannot really be distributed, although that	10:29:44
16	content is public. I mean only constituted of	10:29:47
17	public posts and pictures from -- from users of	10:29:52
18	Meta's platforms. So that would not be distributed.	10:29:57
19	There's also the fact that it's not really	10:30:05
20	useful for that data to be distributed in the first	10:30:09
21	place because the reason for open source --	10:30:12
22	(Stenographer clarification.)	10:30:20

11/21/2024

Richard Kadrey, et al. v. Meta Platforms, Inc.

Yann LeCun, Ph.D.

Highly Confidential - Attorneys' Eyes Only - Under the Protective Order

Page 26

1 THE WITNESS: -- open sourcing a foundation 10:30:22

2 model, such as LLaMA, is -- is to allow people to 10:30:25

3 use this model without having to train it from 10:30:30

4 scratch which is insanely expensive. So giving the 10:30:32

5 data would not help people particularly. 10:30:41

6 And there's a lot of work that goes into 10:30:45

7 curating a training dataset. Basically filtering it 10:30:52

8 in such a way that it only contains high-quality 10:30:59

9 data, and that is not viewed as something that 10:31:03

10 should be publicly revealed for various reasons. 10:31:09

11 Q. (By Ms. Geman) So my question is: Do you 10:31:13

12 agree in Meta's decision to fail to disclose the 10:31:15

13 training sources for LLaMA 2 and 3? 10:31:18

14 A. So that was a decision by our legal 10:31:21

15 department and I have no input. 10:31:24

16 Q. But you're the chief scientist and a very 10:31:28

17 well-regarded scholar and practitioner of AI, so my 10:31:30

18 question is: Do you agree with that decision? 10:31:36

19 MR. WEINSTEIN: Okay. 10:31:38

20 You can state your personal opinion. Her 10:31:38

21 question isn't asking for it, but don't reveal the 10:31:40

22 communications you had with counsel. 10:31:43

11/21/2024

Richard Kadrey, et al. v. Meta Platforms, Inc.

Yann LeCun, Ph.D.

Highly Confidential - Attorneys' Eyes Only - Under the Protective Order

Page 129

1 plural or not, and so typically a token represents 12:41:44

2 about three-quarter of a word on average. 12:41:47

3 And then there are -- 12:41:50

4 (Stenographer clarification.) 12:41:55

5 THE WITNESS: -- agglutinative languages 12:41:55

6 such as German where a word can be constructed by 12:41:57

7 sticking together many words, right, so you have to 12:42:02

8 break that word into individual components. So 12:42:04

9 that's what a token is, really; it's kind of a unit 12:42:08

10 of language that could be smaller than the word. 12:42:12

11 Q. (By Ms. Geman) So let's just say that was 12:42:16

12 the only sentence the LLM trained on it, it would 12:42:17

13 have a set of tokens for that sentence, right -- 12:42:21

14 A. Yeah. 12:42:23

15 Q. -- "a sheep, a wolf and a cabbage lived on 12:42:24

16 an island." 12:42:28

17 Would -- if you were to give -- and every 12:42:29

18 single time it had -- every single time it trained 12:42:33

19 on that sentence, let's say, it would have that same 12:42:36

20 group of tokens; right? Okay. 12:42:39

21 And if you were to sort of show the LLM 12:42:43

22 that group of tokens, it could spit out that 12:42:47

11/21/2024

Richard Kadrey, et al. v. Meta Platforms, Inc.

Yann LeCun, Ph.D.

Highly Confidential - Attorneys' Eyes Only - Under the Protective Order

Page 130

1 sentence; correct? 12:42:51

2 MR. WEINSTEIN: Object to form. 12:42:54

3 THE WITNESS: That's a tricky -- tricky 12:42:57

4 answer to that question. 12:43:02

5 Q. (By Ms. Geman) That sentence -- 12:43:07

6 A. No, I can answer the question, but -- 12:43:07

7 Q. Okay. 12:43:08

8 A. Okay. 12:43:10

9 Q. Sorry, we're noticing it was an hour and 12:43:12

10 counsel was looking at me, so that's what were 12:43:14

11 noticing. 12:43:15

12 But please -- please go ahead and then we 12:43:15

13 can do our lunch break. 12:43:17

14 A. Okay. 12:43:19

15 The same way that if I give you a short 12:43:19

16 poem and I keep repeating it to you multiple times, 12:43:22

17 you will learn that poem by -- by heart. But you're 12:43:28

18 going to be able to regurgitate that poem, if that's 12:43:33

19 the only thing you hear for two days; right? 12:43:36

20 But if the space of those two days you read 12:43:39

21 two novels, you're not going to be able to 12:43:42

22 regurgitate all the words in those novels. 12:43:46

11/21/2024

Richard Kadrey, et al. v. Meta Platforms, Inc.

Yann LeCun, Ph.D.

Highly Confidential - Attorneys' Eyes Only - Under the Protective Order

Page 131

1 It's exactly the same thing for LLMs. If 12:43:48
2 you train them on a small amount of data and it's a 12:43:50
3 big LLM with a big capacity, memory capacity, and 12:43:54
4 you train it -- you train it on a small amount of 12:43:56
5 text, it's probably going to be able to regurgitate 12:43:58
6 that, the text you've been training it on, with some 12:44:02
7 degree of fidelity. 12:44:05

8 But if you train it on a lot of data, it 12:44:08
9 just does not have the memory capacity to -- to 12:44:11
10 store the content and -- and therefore it's not 12:44:14
11 going to be able to regurgitate any of those texts 12:44:17
12 to any significant extent unless within the training 12:44:21
13 set, a particular short segment of text appears 12:44:26
14 thousands of times. 12:44:30

15 So a famous quote, for example, may appear 12:44:31
16 thousands of times in a training set and the LLM 12:44:35
17 obviously would be able to regurgitate that quote 12:44:38
18 because, you know, it appears so many times and it's 12:44:41
19 probably part of common language if it appears so 12:44:44
20 many times, but it's not going to be able to 12:44:47
21 regurgitate the whole novel or even a chapter or 12:44:50
22 even a paragraph probably. 12:44:53

11/21/2024

Richard Kadrey, et al. v. Meta Platforms, Inc.

Yann LeCun, Ph.D.

Highly Confidential - Attorneys' Eyes Only - Under the Protective Order

Page 132

1 Q. If a large LLM trained on only one book, it 12:44:55
2 could regurgitate -- it could regurgitate that book; 12:44:59
3 isn't that true? 12:45:02

4 A. The answer to that question depends on the 12:45:07
5 length of the book and the size of the LLM. So if 12:45:09
6 you train a small language model on a single book, 12:45:13
7 it's not going to be able to regurgitate it. 12:45:16

8 Q. LLaMA could; correct? 12:45:19

9 A. Well, there's sizes of LLaMA. LLaMA 3 12:45:21
10 has -- 12:45:25

11 (Stenographer clarification.) 12:45:27

12 THE WITNESS: -- several versions of 12:45:27
13 various sizes. The smallest one is 1 billion 12:45:28
14 parameters. And then there is a 7 billion, 70 12:45:33
15 billion and a 405 billion. 12:45:34

16 Q. (By Ms. Geman) Are you -- are you 12:45:36
17 disputing that LLaMA 3 could regurgitate a book if 12:45:39
18 that's all it was trained on? 12:45:44

19 A. If it was trained on a single book, not too 12:45:45
20 long, it might be able to store the entire book in 12:45:50
21 its memory and regurgitate it relatively accurately, 12:45:55
22 but nobody would want an LLM of this type. It would 12:46:00

11/21/2024

Richard Kadrey, et al. v. Meta Platforms, Inc.

Yann LeCun, Ph.D.

Highly Confidential - Attorneys' Eyes Only - Under the Protective Order

Page 133

1 be an extremely inefficient way of storing a book 12:46:03

2 and inaccurate and also expensive. 12:46:08

3 MS. GEMAN: Okay. 12:46:11

4 We can take a break. 12:46:11

5 THE WITNESS: Okay. 12:46:13

6 THE VIDEOGRAPHER: We're going to go off 12:46:14

7 the record. The time is 12:46 p.m. 12:46:15

8 (Lunch recess 12:46 p.m. to 1:42 p.m.) 12:46:26

9 THE VIDEOGRAPHER: We are back on the 12:46:26

10 record. The time is 1:42 p.m. 13:42:19

11 Please proceed. 13:42:21

12 Q. (By Ms. Geman) Hi, Yann. Did you have a 13:42:25

13 good lunch? 13:42:26

14 A. Same as yours, I'm afraid. 13:42:27

15 Q. Okay. 13:42:34

16 I notice you didn't answer if it was a good 13:42:34

17 lunch. 13:42:37

18 What is memorization -- what is LLM 13:42:37

19 memorization? 13:42:42

20 A. So it's something that does not really 13:42:43

21 exist in the sense that -- I think as we alluded to 13:42:44

22 earlier, when you train an LLM on a large 13:42:47

11/21/2024

Richard Kadrey, et al. v. Meta Platforms, Inc.
Highly Confidential - Attorneys' Eyes Only - Under the Protective Order

Yann LeCun, Ph.D.

Page 323

1 I, HEATHER J. BAUTISTA, CSR No. 11600,
2 Certified Shorthand Reporter, certify:

3 That the foregoing proceedings were taken
4 before me at the time and place therein set forth,
at which time the witness declared under penalty of
5 perjury; that the testimony of the witness and all
objections made at the time of the examination were
6 recorded stenographically by me and were thereafter
transcribed under my direction and supervision; that
7 the foregoing is a full, true, and correct
transcript of my shorthand notes so taken and of the
8 testimony so given;

9 (XX) Reading and signing was not requested/offered.

10 I further certify that I am not financially
11 interested in the action, and I am not a relative or
12 employee of any attorney of the parties, nor of any
13 of the parties.

14 I declare under penalty of perjury under the
15 laws of California that the foregoing is true and
16 correct. Dated: November 26, 2024

17

18

19

20

21

22



HEATHER J. BAUTISTA, CSR, CRR, RPR, CLR